# **DARIO Web Server**

A free web server for the analysis of short RNAs from high throughput sequencing data

throughput sequencing data. • fast and coordinates the fast and coordinates of the sequencing data.	home I new job	
<ul> <li>no dependencies to any sequencing platform or management</li> </ul>	Transformer and the	
differential expression: direct comparison of exercise	(M. 1570)	
upload of own annotations     upload of own annotations		
prediction of new putative ncRNAs using a machine	Gir all	
CLICK MEDICAL Validation of predicted ncRNAs using RNAz		1
ALIGA HERE TO START A NEW ANALYSIS JOB	The second se	
Read about the following sections:	Teacher and the second se	
1. Introduction	Ora office	
2. Analysis and Prediction		
4. Output		
S. EAQ		
Introduction		
	100	

## SEQUENCING EXPERIMENT AND ITS RESULTS



Sequence small RNAs using a sequencing machine of your choice.

The result will be a list of sequenced RNA molecules (reads) in a sequencer specific output file (typically in fastq format).

## MAP READS TO A REFERENCE GENOME



Freely choose one of your favorite mapping tools to map your data to a reference genome. For example:

- ▶ segemehl
- ► BWA
- Bowtie
- ► SOAP
- •

@HD	VN:1.0											
@SQ	SN:chrI LN:15	072421										
@SQ	SN:chrII LN:15	279323										
@SQ	SN:chrIII	LN:1378	3681									
@SQ	SN:chrIV LN:17	493785										
@SQ	SN:chrV LN:20	919568										
@SQ	SN:chrX LN:17	718854										
0PG	ID:segemehl	VN:0.9.	4-\$Rev: 162	\$ (\$Date: 2010	-10-15 12:48	:37 +0200	(Fri, 15	Oct 2010) \$	)			
GPL9269_	GSM427346_GSE1	7153_102273	10 ch	ITI 15070506	6255 31M	M1D1M *	0	0	ATTCTTAGTTGGTTGAGCGAT	NM:i:4	MD:Z	XN:i:1
GPL9269	GSM427346_GSE1	7153_102273	10 ch	ITI 15063309	9255 31M	11D1M *	0	0	ATTCTTAGTTGGTTGAGCGAT	NM:i:4	MD:Z	XN:i:1
GPL9269_	GSM427346_GSE1	7153_38401	50 ch	IrIV 3233310	255 21M	1 *	0	0	TGAGATCGTTCAGTACGGCAA	NM:i:0	MD:Z:21	XN:i:1
GPL9269	GSM427346_GSE1	7153_38401	50 ch	IrIV 3233310	255 21M	1 *	0	0	TGAGATCGTTCAGTACGGCAA	NM:i:0	MD:Z:21	XN:i:1
GPL9269	GSM427346 GSE1	7153 384015	50 ch	IrIV 3233310	255 21M	1 *	0	0	TGAGATCGTTCAGTACGGCAA	NM:i:0	MD:Z:21	XN:i:1
	_											

example output in sam format

## PREPARE YOUR MAPPED READS FOR THE DARIO UPLOAD



**Automatically convert** the output to the needed bed format using map2bed.pl. Optionally you can create the bed file yourself.

map2bed.p1 will not only create a bed formatted file, but also merge mapped reads to tags and zip the output file. This will minimize the file size, resulting in a short upload time.

(Note: It might be necessary to explicitly select 'SAM format' as output format when running your mapping tool.)

## UPLOAD YOUR DATA TO THE DARIO WEBSERVER



Open the DARIO WebServer http://dario.bioinf.uni-leipzig.de

- Click on CLICK HERE TO START A NEW ANALYSIS JOB
- Choose your reference species
- Choose upload file
- Optionally: choose a list of your own loci of interest
- Optionally: specify an e-mail address
- Click Proceed

### WAIT A MINUTE...

The job is automatically added to a queue and starts as soon as possible. This page reloads every 30 sec and opens the result page when the job is done.



If you specified your e-mail address, you will get an e-mail with a link to the result page, as soon as the job is finished.

## ...DONE

After ~10 to 25 minutes running time the result page is generated, containing the following sections:

- Summary
- Quality Control
- Analysis
- Prediction
- User Annotation
- Download



### SUMMARY

#### Summary

Click <u>here</u> to bookmark this page for reference. This result will be available online at least 2 weeks from your jobs finish time.

You will find all downloadable files at the end of this page.

Job received at	2011-01-24 17:40:51				
Job finished at	2011-01-24 17:50:04				
Uploaded file name	GSM450600.bed.gz				
# of uploaded mapping loci	655,276				
Total # of reads	6,632,854				
Total # of tags	351,359				

The **Summary** contains some basic information about the job, e.g.:

- The date and time you uploaded your file
- The number of mapped loci you uploaded
- The number of reads and the number of tags\*

\*A tag is defined as a RNA sequence that occurs at least once in a set of sequencer reads. Thus a tag typically corresponds to several identical reads.

## **QUALITY CONTROL**



The **Quality Control** gives a first impression on how good your experiment performed.

## QUALITY CONTROL READ LENGTH DISTRIBUTION



#### **Read Length Distribution**

**Quality Control** 

Here, you can easily assess if the short read sequencing experiment performed properly. When using a short RNA preparation protocol, one would expect a peak at the length of 22.

## QUALITY CONTROL MULTIPLE MAPPINGS DISTRIBUTION



Most of the reads should map to a unique position within the reference genome. If this plot gives any doubt to this statement, you should have a deeper look on the quality values of your experiment.

## QUALITY CONTROL MAPPING REGIONS WITHIN THE GENOME



#### **Mapping Region**

This plot shows the fraction of reads mapping to exons, introns, ncRNAs or notannotated regions. The majority of the reads should map to annotated ncRNAs, if a short RNA protocol was used.

### **QUALITY CONTROL** READS MAPPING TO NON-CODING RNA LOCI

Quality Control



#### ncRNA Overlap

This picture gives an overview of the overall ncRNA expression in the experiment. Typically, miRNAs are the majority.

#### ANALYSIS

#### Analysis

This table summarizes the quantification of of the different ncRNA classes. Click on "View List" for detailed information for individual ncRNA species.

ncRNA Class	Reads	Reads (normalized)	# of Genes	Table
miRNA	1,141,865	816,175.5	661	View List
snoRNA_CD	50,991	29,031.09	181	View Lis
snoRNA_HACA	2,125	2,019.533	85	View Lis
tRNA	433,796	72,332.05	537	View Lis
scRNA	286,672	17,815.27	900	View Lis
snRNA	294,027	52,661.4	712	View Lis
rRNA	92,582	24,588.18	389	View Lis
snoRNA_scaRna	16,770	16,752.33	19	View Lis
misc_RNA	105	95.00559	8	View Lis

The tables in the **analysis section** are itemized by the types of ncRNAs.

## ANALYSIS EXAMPLE: MICRORNAS

ncRNA Class		ds (normalized)							
niRNA	1,141,865 816	,175.5	661 <u>V</u>	liew List	>				
noRNA_CD	xpression	of ncRNA	A: miRNA						<u>Go Ba</u>
n s	niRNA expressio ort order.	n sorted with	location. Clic	k the he	ader to sort the ta	ble with res	pect to an	nother column. Click Reads (normalized)	twice to reverse
	chr10	100144965	100145054	-	hsa-mir-1287	2.44e-01	30	30.00	View at UCSC
	chr10	103351164	103351244	+	hsa-mir-3158-1	1.90e-01	42	21.00	View at UCSC
	chr10	103351164	103351244	-	hsa-mir-3158-2	1.90e-01	42	21.00	View at UCSC
	chr10	104186259	104186331	+	hsa-mir-146b	1.20e+01	1.203	1.202.00	View at UCSC

In the expression lists contains:

- the ncRNA loci (chromosomal position)
- ▶ the ncRNA ID
- ▶ the number of reads overlap with the loci
- the number of overlapping reads, normalized for multiple mappings
- ▶ the RPM (Reads Per Million) normalized expression for each expressed ncRNA

## ANALYSIS EXAMPLE: MICRORNAS



Click *View at UCSC* to be forwarded to the UCSC Genome Browser showing the region of the ncRNA. Take a look on the **expression pattern**.

### PREDICTION



The upper part gives information about the sensitivity of the predictions on the uploaded dataset. Below the predicted loci are shown.

## PREDICTION THE IDEA

Short RNA reads map to ncRNAs forming specific patterns (Langenberger et al.).

This is best known for miRNA, where the reads form high stacks right above the annotated mature and the mature\* loci.



read pattern for hsa-mir-204 in the UCSC genome Browser



read pattern on the folded microRNA

## PREDICTION THE IDEA

Other classes of ncRNAs show different expression patterns:





## PREDICTION THE IDEA

Different features are used to separate the classes, using a RandomForest classifier:



## PREDICTION THE OUTPUT

New ncRNAs candidates are identified in the uploaded data and ranked by their scores.

Predictions of ncRNA: miRNA										
Sorted miRNA	predictions of	descending w	ith score	. Click the h	eader to	o sort the ta	ble with	respect to an	other colum	n. Click twice
o reverse sort	order.									
Chromosome	Start Loci	End Loci	Strand	ID	Score	RPM	Reads	Reads (normalized)	RNAz Validation	Visualization
chr11	67984196	67984253	-	miRNA_94	1	1.62e-01	14	14.00		UCSC
chr14	65007583	65007643		miRNA_1	1	1.27e+00	211	106.00		UCSC
chr11	62091059	62091114	-	miRNA_76	0.99	1.83e-01	14	14.00		UCSC

Futhermore, we overlap the predicted loci with RNAz screens. RNAz uses conservation and secondary structure properties to predict regions likely to form functional RNA structures.

## PREDICTION THE EXAMPLES

Currently unknown microRNA candidates predicted by DARIO:



### **USER ANNOTATION**

#### **User Annotation**

You have provided additional annotation (ncRNAs.bed). The following table shows expression of those.

Class Identifier	# of Genes	Table
snoRNA_CD	181	View List
rRNA	389	<u>View List</u>
snRNA	712	View List
tRNA	537	<u>View List</u>
scRNA	900	<u>View List</u>
snoRNA_HACA	85	View List
miRNA	661	View List
snoRNA_scaRna	19	View List
misc_RNA	8	View List

The output is in the same format as the analysis with the user defined annotations instead of the known ncRNA loci.

If you have not uploaded an annotation file, there will be no User Annotation section on the result page.

### DOWNLOAD

#### Download

The following results are available in BED-file format:

- ncRNA Expressions BED
- Predictions BED
- User Annotation Expression BED

Of course all the expression data and predicted candidates DARIO calculated are downloadable in bed format.

It is possible to use the predicted ncRNA candidates as **User Annotation** in upcoming DARIO runs.

## We hope that you will enjoy working with DARIO!

Dario is a cooperative project of

David Langenberger and Mario Fasold

